

U.S. Patent Application based on PCT/EP98/05793

Summary of DE 195 37 010 A1

DE 195 37 010 A1 concerns a learning method and device for simulating a dynamic process by simultaneous learning of at least two time series each of which representing different process parameters. For each parameter of the process, a particular learning component is provided which bases on process data of the past. An optimum learning result is obtained by a decorrelation technique.

DE 195 37 010 A1 represents technological background with regard to the use of neural networks. It does not disclose a method for detecting the modes of a dynamic system with a drift segmentation model as claimed in the above U.S. patent application.

THIS PAGE BLANK (USPTO)

⑬ BUNDESREPUBLIK
DEUTSCHLAND



DEUTSCHES
PATENTAMT

⑫ Offenlegungsschrift
⑩ DE 195 37 010 A 1

⑥ Int. Cl.⁸:
G 05 B 13/02
G 06 F 15/18

⑳ Aktenzeichen: 195 37 010.4
㉑ Anmeldetag: 4. 10. 95
㉒ Offenlegungstag: 10. 4. 97

DE 195 37 010 A 1

㉗ Anmelder:
Siemens AG, 80333 München, DE

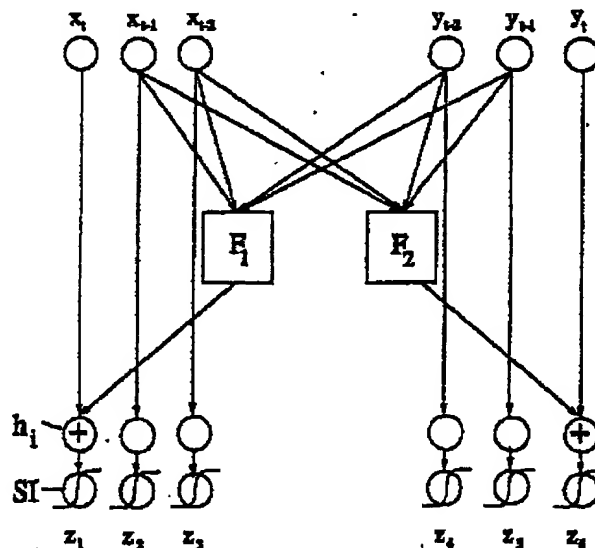
㉘ Erfinder:
Storck, Jan, 82110 Germering, DE; Deco, Gustavo,
Dr., 80636 München, DE

㉙ Entgegenhaltungen:
US 53 98 415
US 51 59 860
US-Z.: Physical Review E, Vol. 51, No. 3, March 1995,
S. 1780-1790;
GB-Z.: Network: Computation in Neural Systems, 5,
1994, S. 565-581;
DE-Z.: atp Automatisierungstechnische Praxis, 37,
1995, 4, S. 55-61;

Prüfungsantrag gem. § 44 PatG ist gestellt

㉚ Lernverfahren und -anordnung zur Nachbildung eines dynamischen Prozesses durch gemeinsames Erlernen von mindestens zwei Zeitreihen, welche jeweils verschiedene Prozeßobservablen darstellen

㉛ Lernverfahren und -anordnung zur Nachbildung eines dynamischen Prozesses durch gemeinsames Erlernen von mindestens zwei Zeitreihen, welche jeweils verschiedene Prozeßobservablen darstellen.
Mit der Erfindung wird eine neuartige Anordnung und ein neuartiges Lernverfahren zur Nachbildung komplexer Prozesse angegeben. Für jede Observable des Prozesses wird eine eigene lernfähige Komponente bereitgestellt, der lediglich Vergangenheitswerte der verwendeten Beobachtungsgrößen zugeführt werden. Dieser Vorgang wird durch die gewählte Anordnung unterstützt. Ein optimales Lernergebnis wird erzielt, indem beim Training die jeweilige Gegenwarts-komponente von ihren Vergangenheitswerten und denen der anderen Beobachtungsgrößen optimal dekorreliert wird. Bei der Durchführung des Verfahrens wird eine besonders günstige Kostenfunktion angegeben. Auf der Architekturseite wird dem Verfahren dadurch Rechnung getragen, daß die Gegenwartswerte der Zeitreihen direkt durchgeschleift werden und lediglich additiv mit der Ausgangsgröße aus der Lernkomponente verknüpft werden.



DE 195 37 010 A 1

Beschreibung

Lernverfahren und -anordnung zur Nachbildung eines dynamischen Prozesses durch gemeinsames Erlernen von mindestens zwei Zeitreihen, welche jeweils verschiedene Prozeßobservable darstellen.

Die Erfindung bezieht sich auf ein neuartiges Lernverfahren und eine vorteilhafte Anordnung zur Durchführung dieses Lernverfahrens zur Nachbildung technischer oder biologischer Prozesse.

Zur Nachahmung komplexer technischer Systeme werden häufig lernfähige Komponenten eingesetzt, um die Prozesse oder Systeme nachbilden zu können. Diesen Systemen ist dabei zueigen, daß sie selbsttätig die Prozeßeigenschaften erlernen können und sich an das Verhalten des nachzubildenden Prozesses anpassen. Insbesondere werden solche Systeme für Prozesse eingesetzt, welche in hohem Maße nicht deterministisch sind, oder die im hohen Grad stochastisch verlaufen. Häufig werden für Steuer- und Regelprobleme in diesem Zusammenhang neuronale Netze oder Fuzzy-Regler eingesetzt.

Bei bisher gängigen Trainingsverfahren für beispielsweise neuronale Netze, werden dem neuronalen Netz Eingangszeitreihen zugeführt und die ausgegebenen Werte des Netzes mit den Eingangswerten verglichen. Der Lernerfolg wird daran gemessen, inwieweit sich die Ausgangswerte den Eingangswerten annähern. Durch gängige Methoden werden die Gewichte an den einzelnen Neuronen eines neuronalen Netzes verändert werden um eine Anpassung, also ein Training des Netzes durchführen zu können. Weitere Lernverfahren sind derzeit nicht bekannt.

Die der Erfindung zugrundeliegende Aufgabe besteht darin, eine Lernanordnung und ein Verfahren anzugeben, womit mehrere verschiedene Observablen eines Prozesses gemeinsam zur Bestimmung einer Ausgangsgröße dieses Lernverfahrens, bzw. dieser Lernanordnung beitragen. Insbesondere soll durch das erfindungsgemäße Verfahren sichergestellt werden, daß nicht eine Ausgangsgröße selbst zur Messung des Lernerfolges herangezogen wird.

Diese Aufgabe wird für das Lernverfahren gemäß den Merkmalen des Patentanspruchs 1 und für die Lernanordnung gemäß den Merkmalen des Patentanspruchs 6 gelöst.

Weiterbildungen der Erfindung ergeben sich aus den abhängigen Ansprüchen.

Ein besonderer Vorteil des erfindungsgemäßen Verfahrens besteht darin, daß zur Bildung einer Gegenwartskomponente alle Vergangenheitskomponenten von Zeitreihen der verschiedensten Observablen herangezogen werden. Besonders vorteilhaft wird durch die optimale Dekorrelation der Gegenwartswerte von allen Vergangenheitswerten sichergestellt, daß der maximal mögliche Lernerfolg eingestellt werden kann.

Um den Rechenaufwand beim erfindungsgemäßen Verfahren und bei der Anordnung vereinfachen zu können werden lediglich die Gegenwartskomponenten durch die Funktionsapproximatoren verändert und die Vergangenheitskomponenten im wesentlichen unverändert an die Ausgänge weitergegeben.

Vorteilhaft werden beim erfindungsgemäßen Verfahren zur einfacheren Weiterverarbeitung und Normierung die auszugehenden Werte mit einer zwischen 0 und 1 beschränkten differenzierbaren Funktion, beispielsweise einer sigmoiden Funktion, bearbeitet.

Besonders vorteilhaft können nach dem erfindungsgemäßen Verfahren Observable danach ausgewählt werden, inwieweit sie nützliche Informationen zum Lernprozeß des jeweiligen Funktionsapproximators beitragen. Ein Maß für diese Nützlichkeit einer solchen Observablen ist das Korrelationsmaß, das zwischen ihr und den anderen Observablen gebildet werden kann. Je weiter diese Observable dekorrelierbar ist, desto nützlicher ist sie für den Lernprozeß des erfindungsgemäßen Verfahrens und einer erfindungsgemäßen Anordnung.

Besonders vorteilhaft wird das erfindungsgemäße Verfahren mit der angegebenen Kostenfunktion durchgeführt, da sie sowohl das Infomax-Prinzip beinhaltet als auch die Korrelation bewertet. Mit dem Infomax-Prinzip wird in diesem Zusammenhang sichergestellt, daß ein Maximum an Information von den Eingängen des Verfahrens, bzw. der Anordnung an die Ausgänge weitergeleitet wird.

Besonders vorteilhaft zur Durchführung des erfindungsgemäßen Verfahrens eignet sich eine Lernanordnung, welche für jede Observable Funktionsapproximationsmittel zur Verfügung stellt. Dadurch, daß diesen Funktionsapproximationsmitteln lediglich die Vergangenheitswerte aller Observablen zugeführt werden, wird schon anordnungsseitig sichergestellt, daß die Gegenwartswerte und Vergangenheitswerte dekorreliert werden können.

Besonders vorteilhaft wird ein solcher Funktionsapproximator in Form eines neuronalen Netzes realisiert, da diese weitestgehend untersucht sind und in beliebiger Vielfalt auch als Emulationsprogramme zur Verfügung stehen.

Im folgenden wird die Erfindung anhand von Figuren weiter erläutert.

Fig. 1 gibt ein Beispiel einer erfindungsgemäßen Anordnung an.

Fig. 2 gibt ein Beispiel für einen technischen Prozeß.

Fig. 3 zeigt Beispiele der Auswirkungen des erfindungsgemäßen Verfahrens nach Anwendung auf den Prozeß in Fig. 2.

In Fig. 1 ist ein Beispiel einer erfindungsgemäßen Lernanordnung dargestellt. Ein vorrangiges Ziel der erfindungsgemäßen Anordnung und des erfindungsgemäßen Verfahrens besteht in der multivariaten Modellierung von Zeitreihen. Beispielsweise werden die zeitlichen Entwicklungen von Systemgrößen eines dynamischen Systems mit Hilfe eines multivariaten Modells auf unüberwachte Weise gelernt. Eingabewerte des Systems sind beispielsweise die Meßwerte mehrerer Observablen des betrachteten Systems. Erfindungsgemäß wird daraus extrahiert, auf welche Weise ein Zeitreihenwert einer Observablen von der eigenen Vergangenheit und von der Vergangenheit weiterer Observabler abhängt. Resultat der erfindungsgemäßen Vorgehensweise ist eine Dekorrelation zwischen der Gegenwart und der Vergangenheit der betrachteten Zeitreihen.

Korrelationen höherer Ordnung, also sowohl lineare als auch nichtlineare Abhängigkeiten zwischen den gemessenen Observablen können dabei extrahiert werden. Diese Korrelationsanalyse gibt beispielsweise Auf-

schluß darüber, ob weitere Meßgrößen eines Systems gegenüber schon gegebenen Observablen auch tatsächlich neue Information über das betrachtete System liefern. Weiterhin kann nach dem Lernvorgang die extrahierte Abhängigkeit zwischen Gegenwart und Vergangenheit zur Vorhersage durch die der Zeitreihenwerte und somit zukünftiger Systemzustände verwendet werden. Diese Prognose gestaltet sich besonders einfach, denn die Funktionsapproximatoren repräsentieren Abbildungen, nach denen sich die Zeitreihen der Observablen zeitlich fortentwickeln. Besonders vorteilhaft kann man das erfindungsgemäße Verfahren und eine Anordnung zur Durchführung des Verfahrens also dafür verwenden, daß die zeitliche Entwicklung einer ganz bestimmten Systemgröße erlernt wird, in dem gelernt wird, wie diese Größe von der eigenen Vergangenheit, als auch von der zusätzlicher anderer Observablen abhängt. Zum anderen können Abhängigkeiten zwischen den verschiedenen Größen erkannt werden.

Besonders vorteilhaft wird durch das erfindungsgemäße Verfahren und eine Anordnung zu dessen Durchführung die Verbindung von unüberwachtem Lernen und multivariater Zeitreihenanalyse hergestellt. Damit gestaltet sich auf erfindungsgemäße Weise die Simultanmodellierung mehrerer Systemgrößen besonders einfach. Insbesondere weist das erfindungsgemäße Verfahren keine Beschränkung auf lineare oder normal verteilte Abhängigkeiten zwischen den Zeitreihenwerten auf. Weiterhin wird durch das erfindungsgemäße Verfahren eine besonders einfache Kostenfunktion zur Verfügung gestellt, welche bezüglich ihrer Anwendung aber eine große Allgemeinheit aufweist.

Die Vorteile des erfindungsgemäßen Verfahrens bestehen insbesondere darin, daß es fähig ist Korrelationen beliebiger Art und Ordnung zu extrahieren. Weiterhin weist es eine besonders niedrige Einbettungsdimension auf, das heißt weniger vergangene Zeitreihenwerte je verwendeter Observabler, als bei univariater Modellierung sind nötig. Besonders günstig wird durch das erfindungsgemäße Verfahren der negative Einfluß von Meßrauschen vermindert. Weiterhin wird durch das erfindungsgemäße Verfahren alle vorhandene Information optimal genutzt, indem sowohl alle zur Verfügung stehenden Observablen, als auch beliebig viele zeitverzögerte Werte dieser Observablen bei der Modellierung Verwendung finden.

Im Stand der Technik sind die Grundlagen der univariaten Zeitreihenmodellierung mit unüberwachtem Lernen in [D595] angegeben. Beispiele zur Phasenraumrekonstruktion mit zeitverzögerten Variablen geben [SYC91] an. Für überwachte Lernverfahren zur Zeitreihenanalyse sind in [LF87] Beispiele angegeben. Die Grundlagen zur Herleitung der erfindungsgemäß angewandten Kostenfunktion ergeben sich aus [NP94] und dem mathematischen Erklärungsteil. Fig. 1 zeigt das multivariate Modell zur Zeitreihenanalyse am Beispiel zweier Observabler und einer jeweils zweidimensionalen Einbettung (es wird zwei Zeitschritte in die Vergangenheit geschaut). Die Zeitreihe der ersten Observablen ist mit x und die der zweiten Observablen mit y bezeichnet. Die entsprechenden Werte der Zeitreihen werden der erfindungsgemäßen Anordnung an den Eingängen zugeführt. Dabei ist zu beachten, daß das erfindungsgemäße Verfahren und die erfindungsgemäße Anordnung sowohl was die Anzahl der simultan eingespeisten Observablen, als auch was die Höhe der Einbettungsdimensionen in jeder Observablen (Anzahl der zeitlich zurückliegenden Werte), welche nicht für alle Observablen gleich sein müssen, beliebig erweitert werden kann. Es werden beispielsweise Vektoren eingegeben, welche sich aus Elementen der Meßreihen der verwendeten Observablen zusammensetzen. Dieses Prinzip ist als Methode der zeitverzögerten Koordinaten (delay coordinates) oder auch als Takens-Methode bekannt. Die Takens-Methode ist dabei eine Methode, die Trajektorien des Phasenraums, bzw. deren Dynamik in einem Einbettungsraum mittels zeitverzögerter Koordinaten zu rekonstruieren. Die Anzahl der dazu benötigten Werte je Rekonstruktionsvektor ist durch die Einbettungsdimension gegeben, die wiederum von der Dimension des Phasenraums bzw. des Attraktors auf dem sich das System bewegt, bestimmt wird. Im Falle zweier Zeitreihen entsteht der Gesamtvektor also beispielsweise aus zwei zeitlich aufeinanderfolgenden Werten einer x - und einer y -Zeitreihe. Jede einzelne Observable trägt dabei $d + 1$ Komponenten zum Eingabevektor bei, wenn d ihre Einbettungsdimension bezeichnet. Weiterhin steuert je Observable eine relativ zu den anderen Werten neueste Komponente zum Eingabevektor bei, die im folgenden als Gegenwartskomponente oder -wert bezeichnet wird. Die übrigen, weiter zurückliegenden Werte werden im folgenden Vergangenheitskomponenten oder -werte genannt. Wie aus Fig. 1 erkannt werden kann, besteht der Eingabevektor also aus x_t, x_{t-1}, x_{t-2} und y_t, y_{t-1}, y_{t-2} . Dabei bezeichnen x_t und y_t die Gegenwartswerte, während x_{t-1}, x_{t-2} und y_{t-1}, y_{t-2} die Vergangenheitswerte repräsentieren. Die für die Anwendung des erfindungsgemäßen Verfahrens nötige Vielzahl solcher Eingabevektoren (Lern-/Trainingsdaten) erhält man durch schrittweises Durchwandern jeweils gesamter Zeitreihen. Sind beispielsweise die Zeitreihenelemente aufsteigend mit 1, 2, 3, ... numeriert, dann besteht der erste Beitrag dieser Zeitreihe zum Gesamteingabevektor aus den Elementen 1, 2, 3, der zweite Beitrag beispielsweise aus den Elementen 2, 3, 4, der dritte aus 3, 4, 5 usw. Es ist beispielsweise ebenfalls möglich bei Anwendung des erfindungsgemäßen Verfahrens die Sprungweite innerhalb der Zeitreihe größer als Eins zu wählen. Beispielsweise werden alle Eingabewerte, bis auf die jeweils zeitlich neuesten jeder Observablen, das heißt genau die Vergangenheitswerte mit einer beispielsweise sigmoiden Übertragungsfunktion

$$f(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

auf den Bereich zwischen Null und Eins beschränkt, ansonsten aber unverändert ausgegeben. Es kann dafür aber auch jede beliebige andere zwischen 0 und 1 beschränkte differenzierbare Funktion verwendet werden. Die Gegenwartskomponenten werden zu den Funktionswerten von Funktionsapproximatoren F_1, F_2 addiert, die sowohl von den Vergangenheitswerten der jeweils eigenen Zeitreihe, als auch von denjenigen der übrigen Zeitreihen abhängen. Dabei wird durch das erfindungsgemäße Verfahren und die Anordnung sichergestellt, daß kein Zeitreihenwert Einfluß hat auf die von ihm aus gesehen zeitlich zurückliegenden Werte. Besonders die

Kausalität des modellierten Prozesses bleibt damit auch im Modell erhalten. Die Funktionsapproximatoren approximieren die Abbildungsvorschriften, welche den zeitlichen Entwicklungen der Zeitreihen zugrundeliegen. Für jede Zeitreihe gibt es beispielsweise einen solchen Approximator. Hier ist für die x-Zeitreihe in Fig. 1 der Funktionsapproximator mit F_1 und für die y-Zeitreihe der Funktionsapproximator mit F_2 bezeichnet. Beispielsweise kann für jeden dieser Funktionsapproximatoren ein eigenes neuronales Netz verwendet werden. Es sind aber auch durchaus andere lernfähige Komponenten in diesem Zusammenhang denkbar. Nach dem erfindungsgemäßen Verfahren werden die freien Parameter dieser lernfähigen Komponenten, welche die approximierten Funktionen bestimmen, iterativ infolge der Minimierung einer Kostenfunktion immer besser angepaßt. Es liegt also ein Lernvorgang vor. Dieser Lernvorgang wird im folgenden anhand eines Beispiels weiter erläutert.

Nach der Summation der Gegenwarts Komponente mit der Ausgabe des zugehörigen Funktionsapproximators, in Fig. 1 mit einem + gekennzeichnet, erfolgt auch hier beispielsweise die nichtlineare Transformation mit der beispielsweise sigmoiden Übertragungsfunktion, welche nun jedoch einen variablen Parameter α enthält:

$$f(x) = \frac{1}{1 + e^{-\alpha x}} \quad (2)$$

Beim erfindungsgemäßen Verfahren werden die verschiedenen Eingabevektoren beispielsweise als Realisierungen eines stochastischen Prozesses aufgefaßt und produzieren als solche auch eine Wahrscheinlichkeitsverteilung am Ausgang, welche durch die Eingangsverteilung induziert wird. In Fig. 1 sind die Ausgänge mit z bezeichnet. Der Vektor, der die Ausgaben vor der abschließenden nichtlinearen Transformation durch die sigmoide Übertragungsfunktion enthält, heißt im folgenden postsynaptisches Potential. In den Formeln im mathematischen Erklärungsteil und in Fig. 1 wird es mit dem mathematischen Symbol \bar{h} bezeichnet. Seine Komponenten lauten h_i . Diejenigen Komponenten des postsynaptischen Potentials, die von den Vergangenheitswerten abhängen, reproduzieren die Eingangsverteilung. Nur die Verteilung derjenigen Komponenten des postsynaptischen Potentials, welche von den Gegenwarts Komponenten der Zeitreihe herrühren, werden nach dem erfindungsgemäßen Verfahren durch ihren jeweiligen Funktionsapproximator beeinflußt. Falls den zeitlichen Entwicklungen der untersuchten Zeitreihen Abbildungsvorschriften zugrundeliegen, so äußern sich diese in Form statistischer Abhängigkeit zwischen den einzelnen Zeitreihenwerten einer Zeitreihe und auch in Form von Abhängigkeiten zwischen den verschiedenen Zeitreihen. Ein Maß für die statistische Abhängigkeit ist die Redundanz der gemeinsamen (multidimensionalen) Verteilung. Diese Abhängigkeiten liegen auch in der Ausgabeverteilung vor. Eine minimale Redundanz ist erreicht, wenn die Einzelkomponenten voneinander statistisch unabhängig sind. Durch statistische Dekorrelation der zu den Gegenwarts Komponenten gehörenden postsynaptischen Potentiale von den übrigen Komponenten des postsynaptischen Potentials, welche die Eingabeverteilung reproduzieren, kann unter den gegebenen Bedingungen das Minimum in der Ausgaberedundanz erreicht werden. Durch das erfindungsgemäße Verfahren wird so sichergestellt, daß ein maximaler Lernerfolg beim Training erzielt werden kann. Dieses Redundanzminimum ist erreicht, wenn die postsynaptischen Potentiale der Gegenwarts Komponenten konstante Werte liefern, also statistisch unabhängig von den übrigen postsynaptischen Potentialen sind. Die entsprechenden Verteilungen müssen also δ -peak darstellen. Für diesen Fall gilt

$$x_t + F_1(x_{t-1}, x_{t-2}, y_{t-1}, y_{t-2}) = c_1 \quad (3)$$

$$y_t + F_2(x_{t-1}, x_{t-2}, y_{t-1}, y_{t-2}) = c_2 \quad (4)$$

und damit

$$x_t = -F_1(x_{t-1}, x_{t-2}, y_{t-1}, y_{t-2}) + c_1 \quad (5)$$

$$y_t = -F_2(x_{t-1}, x_{t-2}, y_{t-1}, y_{t-2}) + c_2 \quad (6)$$

Die Kostenfunktion für das erfindungsgemäße unüberwachte Lernverfahren muß also zu Redundanzminimierung führen. Denn aus Formel 3 wird deutlich, daß die Funktionsapproximatoren zur Erlangung minimaler Redundanz die funktionalen Abhängigkeiten repräsentieren müssen. Infolge des Dekorrelationsvorganges werden folglich Funktionen erhalten, welche die zeitliche Entwicklung der untersuchten Zeitreihen beschreiben. Im betrachteten Beispiel in Fig. 1 also F_1 und F_2 . Mit diesen Funktionen wird die anschließende Vorhersage zukünftiger Zeitreihenwerte ermöglicht. Zusätzlich muß beispielsweise die im Modell übertragene Information maximiert werden (Linsker's Infomax-Prinzip [Lin88]). Als zu maximierende Funktion, welche beide Anforderungen gleichzeitig erfüllt wird beim erfindungsgemäßen Verfahren vorzugsweise folgender Term verwendet:

$$D(P \parallel \prod_i F_i) = - \int \prod_i dh_i \Psi(\bar{h}) \ln \frac{\Psi(\bar{h})}{\prod_i F_i(h_i)} \quad (7)$$

Dieser Term stellt die Kullback-Leibler-Distanz zwischen multidimensionaler postsynaptischer Potentialverteilung und dem Produkt der Ableitung der Übertragungsfunktionen am Ausgang, beispielsweise als sigmoide Funktion gegeben durch

$$f'(x) = \alpha f(x)(1 - f(x)) \quad (8)$$

dar. Zur Maximierung der Gleichung 7, bzw. Minimierung der Gleichung 9, also sowohl zur Gewichtsadaption beispielsweise der neuronalen Netze, welche die einzelnen Funktionsapproximatoren bilden, als auch für die Optimierung der Parameter α_1 und α_2 der Übertragungsfunktionen für die mit dem Gegenwarts-komponenten korrespondierenden Ausgaben, kann beispielsweise Alopex [UV94], ein Standardoptimierungsverfahren für neuronale Netze verwendet werden. Bei der Implementierung läßt sich als Approximation für das Integral aus Gleichung 7 die Summe

$$\frac{1}{M} \sum_{m=1}^M \ln \frac{\hat{\Psi}(\bar{h}^m)}{\prod_{i=1}^p f'(h_i^m)} \quad (9)$$

verwenden, die dann als Kostenfunktion im erfindungsgemäßen Verfahren minimiert wird. Darin bedeutet p die Anzahl der Ausgabewerte, hier in diesem Beispiel $p = 6$, M die Anzahl der Eingabemuster und \bar{h}^m bzw. h_i^m das multi- bzw. eindimensionale postsynaptische Potential, welches vom m -ten Muster erzeugt wurde. Die multidimensionale Dichte Ψ wird beispielsweise mit Histogrammen durch Boxcounting geschätzt:

$$\hat{\Psi}(\bar{h}^m) = \frac{1}{Ml^p} \cdot ZZ \quad (10)$$

wobei M wieder die Anzahl der Eingabemuster ist, \bar{h}^m das postsynaptische Potential, das vom m -ten Eingabemuster erzeugt wird, und ZZ die Anzahl der Punkte im Würfel bezeichnet, der den Wert \bar{h}^m enthält. Mit l ist darin die Kantenlänge des Würfels benannt. Die sigmoiden Funktionen, welche auf die postsynaptischen Potentiale angewendet werden sind in Fig. 1 am Beispiel von z_1 mit SI bezeichnet. Die Wirkung der Anwendung des erfindungsgemäßen Verfahrens und der erfindungsgemäßen Lernanordnung wird in Fig. 2 und 3 verdeutlicht.

Als technischer Prozeß wird beispielsweise ein Beispiel aus der Strömungsdynamik, das Taylor-Couette-System gezeigt. Das Taylor-Couette-System besteht aus zwei coaxialen Kreiszylindern Z1 und Z2, deren Zwischenraum mit einer Flüssigkeit gefüllt ist. Der innere Zylinder Z1 rotiert um die gemeinsame Achse in Fig. 2 mit GA bezeichnet und verursacht damit ab einer bestimmten Drehzahl, die Rotation ist durch einen Pfeil R symbolisiert, die Bildung stationärer gegensinnig rotierender Taylor-Wirbel. In Fig. 2 sind diese Taylor-Wirbel als KS gekennzeichnet. Der äußere Zylinder ist zur Veranschaulichung des Zusammenhanges hier durchsichtig dargestellt. Bei diesem Beispiel wird von einem Zustand stationärer Taylor-Wirbel mit leicht ausgebildeter Turbulenz ausgegangen. Das Beispiel verdeutlicht die Überlegenheit multivariater Modellierung, hier am Beispiel der Verwendung einer zweiten Zeitreihe, gegenüber univariater Modellierung. Für diesen experimentellen Befund werden zwei Zeitreihen durch Messung axialer Geschwindigkeitskomponenten an den Wirbeln A und B gewonnen. Diese beiden Observablen führen zu zwei verschiedenen Zeitreihen im folgenden ebenfalls mit A bzw. B bezeichnet. Der Ergebnisse des erfindungsgemäßen Verfahrens sind für die zwei verschiedenen Observablen in Fig. 3 untereinander dargestellt. Zur Darstellung der Ergebnisse wurden die Zeitreihen sowohl einzeln, als auch simultan dekorreliert. Die Modellierung mit einer Zeitreihe, univariat bedeutet, daß dem zur jeweiligen Zeitreihe gehörenden Funktionsapproximator nur die Vergangenheitswerte der eigenen Zeitreihe zur Verfügung gestellt wurden. Überkreuzkorrelationen können im univariaten Fall nicht genutzt werden.

Dargestellt sind in Fig. 3 die postsynaptischen Potentiale der Gegenwarts-komponenten der Zeitreihen A (links) und B (rechts) für jedes Eingabemuster. Unter a, das heißt in den obersten beiden Diagrammen werden die Werte vor dem Dekorrelationsvorgang, das heißt bei zufälliger Wahl der Modellparameter in den Funktionsapproximatoren dargestellt. Da, wie zuvor bereits erwähnt wurde, im Idealfall die Funktionen einen δ -peak repräsentieren sollen ist die Blickrichtung auf die Diagramme vorgegeben. Sie ist hier mit P bezeichnet. Es kann erkannt werden, daß unter a sowohl die Zeitreihe A und B sehr weit streuen. Unter b sind die Ergebnisse für univariate Dekorrelation dargestellt. Diese univariate Dekorrelation ist nicht Gegenstand der erfindungsgemäßen Anordnung und des erfindungsgemäßen Lernverfahrens. Sie dient lediglich zur Veranschaulichung des durch die Erfindung gegebenen technischen Fortschritts. Unter c sind letztlich die Ergebnisse für Dekorrelation mit zwei Zeitreihen, also bivariate Dekorrelation dargestellt. Deutlich kann hier erkannt werden, daß aus der Blickrichtung P betrachtet nahezu ein δ -peak vorliegen. Deutlich können auch gegenüber b die schmalere Streubereiche der Kurven erkannt werden. Falls nun als Gedankenbeispiel unter c eine Kurve mit ähnlicher Streubreite vorläge, wie die unter b für die Zeitreihe A, so würde dies bedeuten, daß die zusätzlich zur besseren Dekorrelation von Zeitreihe A gewählte Observable B, aus welcher die Zeitreihe B gebildet wurde, keine zusätzliche Information für das Lernen des Funktionsapproximators von A liefert. Es sollte also vorzugsweise eine andere Observable gewählt werden, welche zu einer Verbesserung des Dekorrelationsergebnisses führt. Die detaillierten Zusammenhänge sind im nun folgenden mathematischen Erläuterungsteil weiter dargestellt.

Mathematischer Erklärungsteil

Im folgenden fassen wir auch das Gesamtmodell als Netz auf und bezeichnen entsprechend Ein- und Ausgabewerte als Neuronen. Falls nichts anderes erwähnt wird, sind alle verwendeten Größen vektoriell zu verstehen.

Jedes einzelne Neuron eines Netzes errechnet aus seiner mehrdimensionalen Eingabe v seine Aktivierung (Ausgabe) in zwei Schritten. Zunächst wird v mit einem Gewichtsvektor w skalarmultipliziert und nach diesem ersten Verarbeitungsschritt entsteht das postsynaptische Potential h :

$$h(v) = \sum_i w_i v_i = w \cdot v. \quad (1)$$

Das postsynaptische Potential h ist also eine deterministische Funktion des Eingangssignals des Neurons. Mit der nichtlinearen Transferfunktion f wird es auf das Ausgangspotential V (Aktivierung des Neurons) abgebildet:

$$V = f(h). \quad (2)$$

Hierbei ist f eine beliebige nichtlineare Funktion, die aber zwischen 0 und 1 beschränkt und invertierbar sein soll. In Betracht kommt z. B. die sigmoide Funktion

$$f(x) = \frac{1}{1 + e^{-\alpha x}} \quad (3)$$

mit der Ableitung

$$f'(x) = \alpha f(x)(1-f(x)). \quad (4)$$

wobei der Parameter α die Steigung und damit den Bereich nahezu linearer Abbildung gegenüber nichtlinearer Übertragung bestimmt.

Wir betrachten nun speziell die Neuronen der Ausgabeschicht. Die Dimension der Ausgabeschicht sei p . Erweitert auf den allgemeinen Fall mehrerer Neuronen sind h und V als vektorielle Größen zu verstehen. Das mehrdimensionale Eingangssignal des Netzes ζ induziert das postsynaptische Potential h mit Verteilung $\Psi(h)$ am Ausgang. Daher ist h eine deterministische Funktion des Zufallsvektors ζ , wobei h beliebige nichtlineare Transformationen enthalten kann. Liegen nämlich eine oder mehrere nichtlineare Schichten zwischen Eingabe- und Ausgabeschicht, dann stellt das Netz einen allgemeinen Funktionsapproximator dar. Derartige Transformationen zwischen Eingabe ζ und postsynaptischem Potential h sind nicht notwendigerweise bijektiv. Es kann also etwas von der Eingangsinformation bei der Übertragung durch das Netz verlorengehen. Unser Ziel ist es nun, die Transinformation $I(\zeta, V)$ zwischen Eingabe und Ausgabe des Netzes zu maximieren, um so eine möglichst verlustfreie Übertragung zu gewährleisten. Da informationstheoretische Größen nur für Zufallsvariablen definiert sind, müssen wir zusätzlich künstliches Rauschen z mit Verteilung $v(z)$ am Ausgangspotential V hinzufügen. Wir erhalten die Aktivierungen der Ausgangsneuronen V als einen zweiten Zufallsvektor

$$V = f(h) + z. \quad (5)$$

wobei f eine invertierbare Transferfunktion mit $0 < f_i < 1$ für alle Komponenten $i = 1, \dots, p$ ist. Für die einzelnen Ausgangsaktivierungen haben wir also

$$V_i = f_i(h_i) + z_i \text{ für } i = 1, \dots, p. \quad (6)$$

Neben den durch die jeweiligen Gewichte vorgegebenen Potentialen h_i können sich auch die Transferfunktionen f_i von Neuron zu Neuron unterscheiden. Aufgrund des lediglich theoretischen Zwecks ist die Wahrscheinlichkeitsverteilung $v(z)$ des additiven Rauschens z hierbei beliebig, wobei z jedoch als unabhängig von h angenommen wird (die z_i 's müssen keine untereinander unabhängigen Zufallsvariablen sein). Die Rauschstärke sei dabei wie folgt definiert:

$$\sum_i (\langle z_i^2 \rangle - \langle z_i \rangle^2) = p\Delta, \quad (7)$$

wobei Δ die Rauschstärke eines einzelnen Ausgabeneurons bezeichnet und $\langle \rangle$ Mittelung über die $v(z_i)$ -Verteilung bedeutet.

Zusätzlich zur Transinformation $I(\zeta, V)$ zwischen Eingabe und Ausgabe betrachten wir nun die Transinformation $I(h, V)$ zwischen dem Potential h und der Ausgabe. Unter der Voraussetzung, daß kein Eingangsrauschen vorhanden ist, sind $I(\zeta, V)$ und $I(h, V)$ gleich. Daher können wir die weitaus handlichere Größe $I(h, V)$ betrachten, um den Informationstransfer des Netzwerkes zu maximieren. Im folgenden wollen wir deshalb einen analytischen Ausdruck für $I(h, V)$ herleiten, der nur von den adaptierbaren Netzparametern abhängt (vgl. [NP94]). Die Transinformation zwischen den Zufallsvektoren h und V ist gegeben durch

$$I(V, h) = \iint \Psi(h) Q(V|h) \ln \frac{\Psi(h) Q(V|h)}{\Psi(h) Q(V)} dh dV. \quad (\overline{8})$$

Hierbei ist $Q(V|h)$ die bedingte Wahrscheinlichkeit von V bei bekanntem h und ergibt sich gemäß (5) zu:

$$Q(V|h) = v(V - f(h)). \quad (9)$$

Als resultierende Ausgangsverteilung erhält man:

$$q(V) = \int \Psi(h) Q(V|h) dh. \quad (10)$$

Aufgrund der Additivität des Rauschens läßt sich die Transinformation I auch als Differenz zwischen den Entropien der Ausgangs- und Rauschverteilungen darstellen:

$$I = H(q) - H(v). \quad (11)$$

Der erste Term in (11) ist die differentielle Entropie der Wahrscheinlichkeitsverteilung q :

$$H(q) = - \int q(V) \ln q(V) dV. \quad (12)$$

Der zweite Term in (11) hängt nur von der Verteilung des Rauschens ab:

$$H(v) = - \int (z) \ln(v(z)) dz. \quad (13)$$

Im Fall, daß v_i ($i=1, \dots, p$) eine Gaußverteilung ist, ist $H(v_i)$ gleich $\frac{1}{2} \ln(2\pi e\Delta)$. Da die Gaußverteilung die größte Entropie unter allen Verteilungen gegebener Varianz hat, gilt

$$H(v_i) \leq \frac{1}{2} \ln(2\pi e\Delta). \quad (\overline{14})$$

Wenn also Δ gegen null geht, streben die Einzelentropien $H(v_i)$ gegen minus unendlich. Es folgt dann, daß damit auch die gemeinsame Entropie gegen minus unendlich geht. Der zweite Term aus (11) strebt also gegen unendlich. Von den beiden Größen aus (11) ist für uns aber lediglich $H(q)$ von Interesse, da sich nur $H(q)$ durch die Adaption von f bzw. der Gewichte beeinflussen läßt. Um die Transinformation I zu maximieren, gilt es also, die Ausgangsentropie $H(q)$ zu maximieren. Für eine gegebene Rauschstärke erzwingt diese Maximierung der Entropie die Bijektivität der Transformation von ζ nach h , was ja genau unser Ziel war. Dies folgt aus der Tatsache, daß Nichtbijektivität eine niedrigere Entropie nach sich zieht. Werden mehrere Eingabewerte auf gleiche Ausgabewerte abgebildet, dann nimmt die Unsicherheit im Ausgabecode und damit auch die Entropie ab. Diese Argumentation gilt allerdings nur, weil die Ausgangstransferfunktionen beschränkt sind. Diese Einschränkung sichert zu, daß die Ausgangsentropie nicht ad infinitum erhöht werden kann, indem der Bildbereich der erzeugten Ausgabe gestreckt wird. Ab einem bestimmten Stadium bleibt dem Netz folglich zu einer weiteren Erhöhung der Entropie lediglich das Mittel der Bijektivität übrig.

Im Limes verschwindenden Rauschens hat die Größe $H(q)$ einen endlichen Grenzwert. Für $\Delta \rightarrow 0$ wird q zu

$$q(V) = \int \Psi(h) \delta(V - f(h)) \prod_j dh_j \quad (\overline{15})$$

$$= \int \Psi(h) \prod_i \delta(V_i - f_i(h_i)) \prod_j dh_j. \quad (\overline{16})$$

Eingesetzt in (12) ergibt sich $H(q)$ zu

$$H(q) = - \int \prod_i dV_i \int \prod_j dh_j \Psi(h) \delta(V_j - f_j(h_j)) \quad (\overline{17})$$

$$\ln \int \prod_k d\tilde{h}_k \Psi(\tilde{h}) \delta(V_k - f_k(\tilde{h}_k)) \quad (\overline{18})$$

$$= - \int \prod_j dh_j \Psi(h) \ln \int \prod_k d\tilde{h}_k \Psi(\tilde{h}) \delta(f_k(h_k) - f_k(\tilde{h}_k)). \quad (\overline{19})$$

Um die restlichen Delta-Integrationen ausführen zu können, machen wir die Substitutionen

$$\bar{y}_i = f_i(\bar{h}_i) \quad (\bar{y}_i \text{ als Funktion von } \bar{h}_i) \quad (\bar{20})$$

5 und

$$y_i = f_i(h_i) \quad (y_i \text{ nimmt den festen Wert } f_i(h_i) \text{ an}) \quad (\bar{21})$$

10 und wir erhalten schließlich

$$H(q) = - \int \prod_j dh_j \Psi(h) \ln \int \prod_k d\bar{y}_k \frac{\Psi(f^{-1}(\bar{y}))}{f'_k(f_k^{-1}(\bar{y}_k))} \delta(y_k - \bar{y}_k) \quad (\bar{22})$$

$$15 = - \int \prod_j dh_j \Psi(h) \ln \frac{\Psi(f_1^{-1}(y_1), \dots, f_p^{-1}(y_p))}{\prod_k f'_k(f_k^{-1}(y_k))} \quad (\bar{23})$$

$$20 = - \int \prod_j dh_j \Psi(h) \ln \frac{\Psi(h)}{\prod_k f'_k(h_k)} \quad (\bar{24})$$

Für die Entropie $H(q)$ und damit für den relevanten Teil der Transinformation I erhalten wir somit den Ausdruck

$$25 \quad H(q) = -D(\Psi \| \prod_k f'_k), \quad (\bar{25})$$

30 wobei

$$D(\Psi \| \prod_k f'_k) \equiv \int \prod_j dh_j \Psi(h) \ln \frac{\Psi(h)}{\prod_k f'_k(h_k)} \quad (\bar{26})$$

35 Da wir $0 < f_i < 1$ für alle $i = 1, \dots, p$ angenommen haben, erfüllt jedes f_i die Voraussetzung einer Wahrscheinlichkeitsverteilung (Integration von $-\infty$ bis $+\infty$ ergibt sich zu eins). Damit kann man dann $D(\Psi \| \prod_k f'_k)$ als Kullback-Leibler-Distanz zwischen der Potentialverteilung Ψ und der Wahrscheinlichkeit auffassen, die durch das Produkt der f_i definiert ist. Ihr Wert ist immer größer oder gleich null, wobei null genau dann angenommen wird, wenn die beiden Verteilungen (bis auf Nullmengen) identisch sind.

Wir halten fest: die Transinformation ist bis auf eine Konstante (gegeben durch die Rauschentropie) gleich minus der Kullback-Leibler-Distanz zwischen der Potentialverteilung und der Produktverteilung, die durch die Ableitungen der Transferfunktionen dargestellt wird. Maximierung der Transinformation ist äquivalent zur Minimierung der Kullback-Leibler-Distanz. Der optimale Fall von $D = 0$ wird genau dann erreicht, wenn

$$\Psi(h) = \prod_i f_i(h_i) \quad (\bar{27})$$

50 gilt. Damit wird außerdem klar: ein faktorieller Code von $\Psi(h)$, d. h.

$$\Psi(h) = \prod_i \Psi_i(h_i), \quad (\bar{28})$$

ermöglicht eine Maximierung der übertragenen Information. Die optimalen Transferfunktionen ergeben sich dann einfach zu

$$60 \quad f'_i(h_i) = \Psi_i(h_i), \text{ für } i = 1, \dots, p \quad (\bar{29})$$

und können für jedes Neuron unabhängig von den anderen adjustiert werden.

Faktorisierung der Verteilung des postsynaptischen Ausgangspotentials ist aber gleichbedeutend mit Redundanzminimierung. Als Ergebnis dieses Abschnitts erhalten wir damit:

65

REDUNDANZMINIMIERUNG



INFORMATIONSMAXIMIERUNG,

(30)

5

unter der Voraussetzung, daß die Transferfunktionen gemäß (29) optimal angepaßt werden.

Einige Bemerkungen: da wir von f_i zunächst nur Invertierbarkeit gefordert haben, käme auch eine streng monoton fallende Funktion mit negativer Ableitung als Transferfunktion in Frage. In den Gleichungen (15) bis (29) wäre dann die allgemeinere Form mit $|f_i(h_i)|$ anstelle von $f_i(h_i)$ zu verwenden und man erhielte als alternative Lösung für (29) $F_i = -\Psi_i$. Wir wollen uns aber im folgenden auf die sigmoide Funktion aus (3) beschränken, so daß wir diesen Fall ausschließen können.

In der Bildverarbeitung ist das Resultat (29) unter dem Namen "Sampling/Histogram Equalization" bekannt. Es besagt, daß maximale Informationsübertragung bei uniformer Ausgangsverteilung — also bei der Verteilung maximaler Entropie — erreicht werden kann.

Physikalisch gesehen läßt sich dieses Ergebnis leicht plausibel machen: eine große Menge an Information wird dann übertragen, wenn das Eingangssignal am Ausgang wieder fein aufgelöst werden kann. Bei Stichproben der empirisch ermittelten Verteilung $\Psi_i(h_i)$ beobachtet man die meisten Stichprobenwerte in der Nähe der h_i -Werte, für die $\Psi_i(h_i)$ groß ist. Um diese gut voneinander trennen zu können, muß dort auch die Steigung der Transferfunktion möglichst groß sein. Verschiedene Ausgangswerte liegen somit weit auseinander und können trotz Rauschens noch unterschieden werden. Eine untere Schranke für die Auflösung ist dabei durch die vom Rauschen bedingte Skalierung am Ausgang gegeben. Die Rauschstärke, unendlich klein, aber ungleich null, setzt also ein Maß für die Trennschärfe der Informationsübertragung.

Nachdem wir im letzten Abschnitt gesehen haben, daß ein faktorieller Code bei entsprechender Wahl der Transferfunktionen maximalen Informationstransfer garantiert, wollen wir nun auch noch die entgegengesetzte Richtung zeigen: Maximierung der Transinformation führt zu einem faktoriellen Code, falls ein solcher existiert. Die Redundanz R im Ausgabe-Code, bedingt durch Korrelationen zwischen den einzelnen Ausgabewerten, ist definiert als

$$R = \sum_i H(q_i) - H(q). \quad (31)$$

30

Für die eindimensionalen Entropien $H(q_i)$ und die multidimensionale Entropie $H(q)$ setzen wir jetzt den im letzten Abschnitt hergeleiteten Ausdruck (26) für die einzelne und für die gemeinsame Entropie ein:

$$R = - \sum_j \int \Psi(h_j) \ln \frac{\Psi(h_j)}{f_j(h_j)} dh + \int \Psi(h) \ln \frac{\Psi(h)}{\prod_i f_i(h_i)} \prod_j dh_j. \quad (32)$$

40

Da die Redundanz R immer nichtnegativ ist, gilt mit

$$D_j = \int \Psi(h_j) \ln \frac{\Psi(h_j)}{f_j(h_j)} dh \quad \text{und} \quad D = \int \Psi(h) \ln \frac{\Psi(h)}{\prod_i f_i(h_i)} \prod_j dh_j$$

45

$$R = - \sum_j D_j + D \geq 0 \quad (33)$$

50

und damit auch

$$D \geq \sum_j D_j. \quad (34)$$

55

Bei den einzelnen Summanden von $\sum_j D_j$ handelt es sich aber lediglich um Kullback-Leibler-Distanzen, so daß auch diese die Bedingung der Nichtnegativität erfüllen. Man erhält schließlich die Ungleichungskette

60

$$D \geq \sum_j D_j \geq 0, \quad (35)$$

d. h.

$$D \rightarrow 0 \Rightarrow \sum_j D_j \rightarrow 0. \quad (36)$$

65

Eine Maximierung der Transinformation I und die damit verbundene Minimierung der Kullback-Leibler-Distanz D Transferfunktionen gegebenen Dichten führt also zwangsläufig zur Minimierung der Redundanz, falls ein faktorieller Code existiert. In unserem speziellen Fall invertierbarer und beschränkter Transferfunktionen, nicht vorhandenem Eingangsrauschen und verschwindend geringem, d. h. infinitesimal kleinem, aber positivem Ausgangsrauschen erhalten wir zusammen mit (30) das Hauptergebnis dieses gesamten Kapitels über Informationsverarbeitung in neuronalen Netzen:

REDUNDANZMINIMIERUNG



INFORMATIONSMAXIMIERUNG,

(37)

unter der Voraussetzung, daß ein faktorieller Code existiert (ist dies nicht der Fall, dann soll die Potentialverteilung wenigstens so weit wie möglich faktorisiert werden). Es ist allerdings zu beachten, daß es genau genommen nur die Informationsmaximierung ist, die sowohl die Parameter für die Transformation T und damit die Potentialverteilung Ψ als auch die Transferfunktionen f_i vorschreibt.

Dieses Ergebnis hat eine fundamentale Bedeutung für unüberwachte Lernverfahren: die Kostenfunktion reduziert sich auf den Infomax-Term, d. h. die Kullback-Leibler-Distanz (26), die das neuronale Netz minimieren soll. Es ist wichtig zu bemerken, daß das Minimum $D = 0$ nur erreicht werden kann, falls die Transformation T und die Transferfunktionen f_i allgemein bzw. flexibel genug sind.

Literatur

- [DS95] Deco, G.; Schürmann, B.: "Learning time series evolution by unsupervised extraction of correlations". — In: Phys. Rev. E 51 (1995), S. 1780—1785.
 [LF87] Lapedes, A.; Farber, R.: Nonlinear signal processing using neural networks: prediction and signal modelling. Technischer Bericht LA-UR-987-2662, Los Alamos National Laboratory, Los Alamos, NM, unveröffentlicht, 1987.
 [Lin88] Linsker, R.: "Self-organization in a perceptual network". — In: IEEE Computer 21 (1988), S. 105—117.
 [NP94] Nadal, J.-P.; Parga, N.: "Non-linear neurons in the low noise limit: a factorial code maximizes information transfer". — In: Network 5 (1994), S. 565—572.
 [SYC91] Sauer, T.; Yorke, J.; Casdagli, M.: "Embedology". — In: J. Stat. Phys. 65 (1991), S. 579—617.
 [UV94] Unnikrishnan, K. P.; Venugopal, K. P.: "Alopex: A correlation-based learning algorithm for feedforward and recurrent neural networks". — In: Neural Computation 6 (1994), S. 469—473.

Patentansprüche

1. Lernverfahren zur Nachbildung eines dynamischen Prozesses durch gemeinsames Erlernen von mindestens zwei Zeitreihen, welche jeweils verschiedene Prozeßobservable darstellen,
 - a) bei dem jede Prozeßobservable durch einen Funktionsapproximator nachgebildet wird,
 - b) bei dem jedem Funktionsapproximator lediglich in der Vergangenheit liegenden Werte aller Zeitreihen zur Verfügung gestellt werden,
 - c) bei dem die einzelnen Werte einer jeweiligen Zeitreihe aufgefaßt werden als mit einer je Wert spezifischen Wahrscheinlichkeitsverteilung auftretende Realisierungen eines stochastischen Prozesses,
 - d) und bei dem zum Training des Funktionsapproximators, der von ihm erzeugte Wert zum jeweiligen Gegenwartswert der Zeitreihe in Form eines Ausgabewertes addiert wird und vom Funktionsapproximator als Ausführungsfunktion eine solche Funktion erzeugt wird, die sicherstellt, daß die Wahrscheinlichkeitsverteilung dieses Ausgabewertes von der Wahrscheinlichkeitsverteilung aller zugeführten Werte optimal dekorreliert ist.
2. Verfahren nach Anspruch 1, bei dem alle Vergangenheitswerte der Zeitreihen identisch ausgegeben werden.
3. Verfahren nach einem der vorangehenden Ansprüche, bei dem auf alle auszugebenden Werte eine differenzierbare Übertragungsfunktion angewendet wird, welche ihnen einen Wert zwischen 0 und 1 zuweist.
4. Verfahren nach einem der vorangehenden Ansprüche, bei dem die Zeitreihe einer bisher nicht verwendeten Prozeßobservable zugeführt wird, falls mit den aktuell verwendeten Zeitreihen keine Dekorrelation möglich ist.
5. Verfahren nach einem der vorangehenden Ansprüche, bei dem zur Einstellung der Ausführungsfunktion am jeweiligen Funktionsapproximator folgende Funktion maximiert wird:

$$D(\Psi || \prod_k f_k) = - \int \prod_j dh_j \Psi(\vec{h}) \ln \frac{\Psi(\vec{h})}{\prod_k f_k(h_k)} \quad (7)$$

mit:

f' : Ableitung der Übertragungsfunktion (8)

\bar{h} : multidimensionales postsynaptisches Potential, bestehend aus allen Vergangenheitswerten und den Summen von Gegenwartswerten mit den Ausgaben der Funktionsapproximatoren

Ψ : multidimensionale Wahrscheinlichkeitsverteilung am Ausgang

D : Kullback-Leibler Distanz

6. Verfahren nach Anspruch 5, bei dem das Integral in Gleichung (7) durch folgenden, als Kostenfunktion zu minimierenden Term angenähert wird:

$$\frac{1}{M} \sum_{m=1}^M \ln \frac{\hat{\Psi}(\bar{h}^m)}{\prod_{i=1}^p f'(h_i^m)} \quad (9)$$

mit:

M : Anzahl der Eingabemuster

\bar{h}^m : multidimensionales postsynaptisches Potential

h_i^m : eindimensionales postsynaptisches Potential

und bei dem folgende Näherung benutzt wird:

$$\hat{\Psi}(\bar{h}^m) = \frac{1}{Ml^p} \cdot ZZ \quad (10)$$

mit:

l : Würfel, der den Wert \bar{h}^m enthält

ZZ : Anzahl der Punkte im Würfel

p : Anzahl der Ausgabewerte

l : Kantenlänge des Würfels

6. Lernanordnung zur Nachbildung eines dynamischen Prozesses durch gemeinsames Erlernen von mindestens zwei Zeitreihen, welche jeweils verschiedene Prozeßobservablen darstellen,

a) bei der mindestens erste und zweite Funktionsapproximationsmittel zur Nachbildung des Zeitverhaltens der jeweiligen Prozeßobservablen vorgesehen sind,

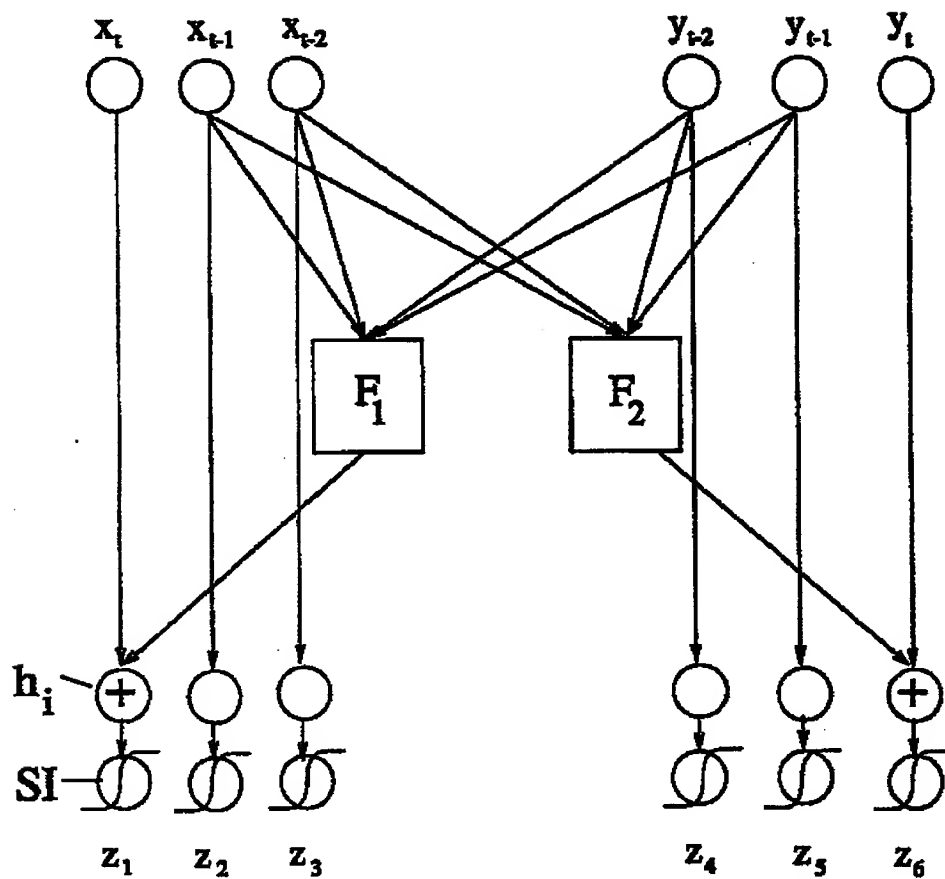
b) bei dem jedem Funktionsapproximationsmittel lediglich in der Vergangenheit liegende Werte aller Zeitreihen zugeführt werden,

c) und bei der im jeweiligen Funktionsapproximationsmittel eine Ausführungsfunktion aus einem der Ansprüche 1—5 ausgeführt wird.

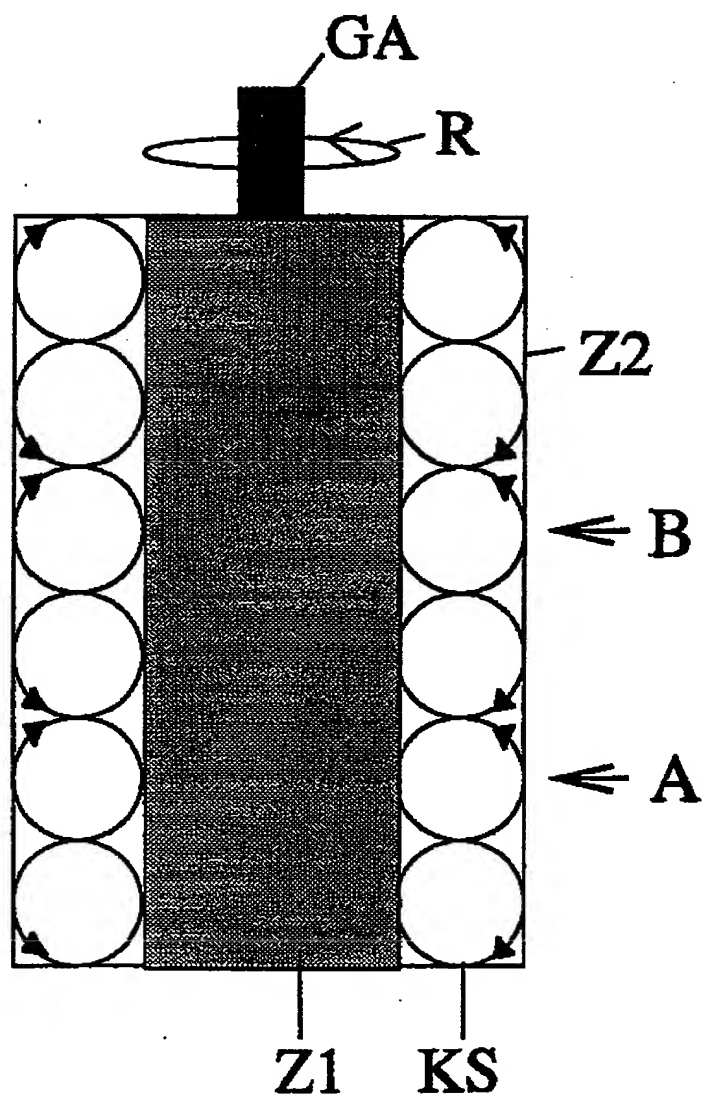
7. Lernanordnung nach Anspruch 6, bei der als Funktionsapproximationsmittel ein neuronales Netz vorgesehen ist.

Hierzu 3 Seite(n) Zeichnungen

Figur 1



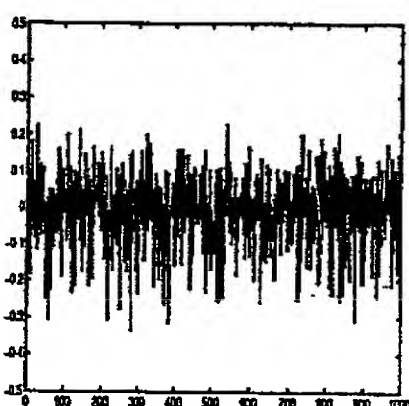
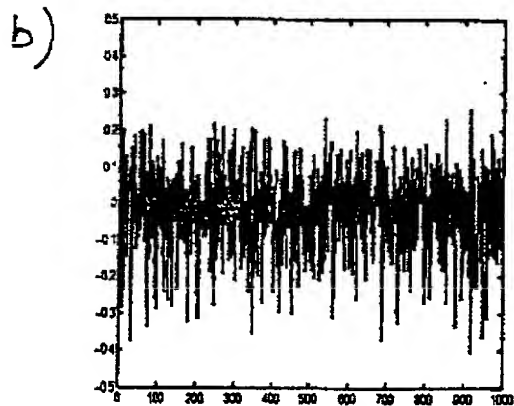
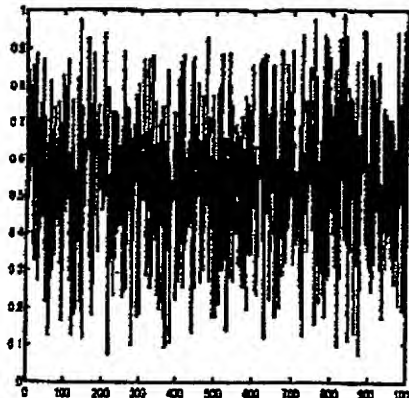
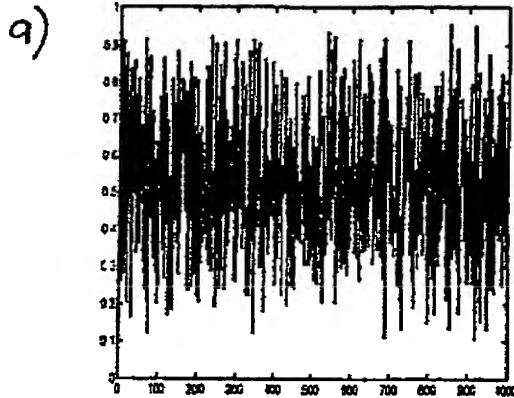
Figur 2



Figur 3

A

B



ρ

